

**System, Method, and Computer Program Product for the
Visualization and Interactive Processing and Analysis of Chemical
Data**

Inventors: Dimitris K. Agrafiotis, Downingtown, PA
Victor S. Lobanov, Exton, PA

Cross Reference to Related Applications

This application is related to and claims priority to U.S. Provisional Application Ser. No. 60/030,187, filed November 4, 1996, titled "Stochastic Algorithms for Maximizing Molecular Diversity," which is herein incorporated by reference in its entirety.

Background of the Invention

Field of the Invention

The present invention is generally directed to displaying and processing data using a computer, and more particularly directed to visualizing and interactively processing chemical compounds using a computer.

Related Art

Currently, research to identify chemical compounds with useful properties (such as paints, finishes, plasticizers, surfactants, scents, drugs, herbicides, pesticides, veterinary products, etc.) often includes the synthesis/acquisition and analysis of large libraries of chemical compounds. More and more, combinatorial chemical libraries are being synthesized/acquired and analyzed to conduct this research.

A combinatorial chemical library is a collection of diverse chemical compounds generated by either chemical synthesis or biological synthesis by combining a number of chemical "building blocks" such as reagents. For example, a linear combinatorial chemical library such as a polypeptide library is formed by combining a set of chemical building blocks called amino acids in every possible way for a given compound length (i.e., the number of amino acids in a polypeptide compound). Millions of chemical compounds theoretically can be synthesized through such combinatorial mixing of chemical building blocks. For example, one commentator has observed that the systematic, combinatorial mixing of 100 interchangeable chemical building blocks results in the theoretical synthesis of 100 million tetrameric compounds or 10 billion pentameric compounds (Gallop *et al.*, "Applications of Combinatorial Technologies to Drug Discovery, Background and Peptide Combinatorial Libraries," J. Med. Chem. 37, 1233-1250 (1994)).

Advanced research in this area often involves the use of directed diversity libraries. A directed diversity library is a large collection of chemical compounds having properties/features/characteristics that match some prescribed properties. The generation, analysis, and processing of directed diversity libraries are described in U.S. Patent Nos. 5,463,564; 5,574,656; and 5,684,711, and pending U.S. Application titled "SYSTEM, METHOD AND COMPUTER PROGRAM PRODUCT FOR IDENTIFYING CHEMICAL COMPOUNDS HAVING DESIRED PROPERTIES," Atty. Docket No. 1503.0200001, all of which are herein incorporated by reference in their entireties.

In conducting such research, it would be very valuable to be able to compare the properties, features, and other identifying characteristics of compounds. For example, suppose that a researcher has identified a compound X that exhibits some useful properties. It would aid the researcher greatly if he could identify similar compounds, since those similar compounds might also exhibit those same useful properties.

It would also help a researcher in his work to be able to easily synthesize compounds, or retrieve compounds from a chemical inventory. Further, it would

greatly aid a researcher to be able to interactively analyze and otherwise process chemical compounds.

Summary of the Invention

5 Briefly stated, the present invention is directed to a system, method, and computer program product for visualizing and interactively analyzing data relating to chemical compounds. The invention operates as follows. A user selects a plurality of compounds to map, and also selects a method for evaluating similarity/dissimilarity between the selected compounds. A non-linear map is generated in accordance with the selected compounds and the selected method. The non-linear map has a point for each of the selected compounds, wherein a distance between any two points is representative of similarity/dissimilarity between the corresponding compounds. A portion of the non-linear map is then displayed. Users are enabled to interactively analyze compounds represented in the non-linear map.

10 Further features and advantages of the present invention, as well as the structure and operation of various embodiments of the present invention, are described in detail below with reference to the accompanying drawings. In the drawings, like reference numbers indicate identical or functionally similar elements. Also, the leftmost digit(s) of the reference numbers identify the drawings in which the associated elements are first introduced.

Brief Description of the Figures

The file of this patent contains at least one drawing executed in color. Copies of this patent with color drawing(s) will be provided by the Patent and Trademark Office upon request and payment of the necessary fee.

20 The present invention will be described with reference to the accompanying drawings, wherein:

FIG. 1 illustrates a block diagram of a computing environment according to an embodiment of the invention;

FIG. 2 is a block diagram of a computer useful for implementing components of the invention;

FIG. 3 is a flowchart representing the operation of the invention in visualizing and interactively processing non-linear maps according to an embodiment of the invention;

FIG. 4 is a flowchart representing the manner in which a non-linear map is generated according to an embodiment of the invention;

FIG. 5 illustrates a structure browser window according to an embodiment of the invention;

FIG. 6 illustrates a compound visualization non-linear map window according to an embodiment of the invention;

FIG. 7 is used to describe a zoom function of the present invention;

FIG. 8 illustrates a dialog used to adjust properties of a set containing one or more compounds;

FIGS. 9 and 10 are used to describe the compound visualization non-linear map window according to an embodiment of the invention;

FIG. 11 is a flowchart illustrating the operation of the invention where a compound visualization non-linear map window is used as a source in an interactive operation;

FIG. 12 is a flowchart illustrating the operation of the invention where a compound visualization non-linear map window is used as a target in an interactive operation;

FIG. 13 conceptually illustrates an interactive operation where a compound visualization non-linear map window is used as a source; and

FIG. 14 conceptually illustrates an interactive operation where a compound visualization non-linear map window is used as a target.

Detailed Description of the Preferred Embodiments

Table of Contents

1.	Overview of the Present Invention	
2.	Structure of the Invention	
5	3. Implementation Embodiment of the Invention	
	4. Overview of Multidimensional Scaling (MDS) and Non-Linear Mapping (NLM)	
	4.1 Procedure Suitable for Relatively Small Data Sets	
	4.2 Procedure Suitable for Large Data Sets	
10	5. Evaluation Properties (Features) and Distance Measures	
	5.1 Evaluation Properties Having Continuous or Discrete Real Values	
	5.2 Distance Measure Where Values of Evaluation Properties Are Continuous or Discrete Real Numbers	
	5.3 Evaluation Properties Having Binary Values	
15	5.4 Distance Measures Where Values of Evaluation Properties Are Binary	
	6. Scaling of Evaluation Properties	
	7. Improvements to Map Generation Process	
	7.1 Pre-Ordering	
	7.2 Localized Refinement	
20	7.3 Incremental Refinement	
	8. Operation of the Present Invention	
	9. User Interface of the Present Invention	
	9.1 Structure Browser	
	9.2 Map Viewer	
25	9.3 Interactivity of the Present Invention	
	9.3.1 Map Viewer as Target	
	9.3.2 Map Viewer as Source	

9.4 Multiple Maps

10. Examples

10.1. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.2. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.3. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.4. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.5. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.6. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.7. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.8. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.9. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.10. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.11. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.12. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.13. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.14. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.15. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.16. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.17. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.18. \mathbb{C}^n is a vector space over \mathbb{C} .
 10.19. \mathbb{R}^n is a vector space over \mathbb{R} .
 10.20. \mathbb{C}^n is a vector space over \mathbb{C} .

1. *Overview of the Present Invention*

The present invention is directed to a computer-based system, method, and/or computer program product for visualizing and analyzing chemical data using interactive multi-dimensional (such as 2- and/or 3-dimensional) non-linear maps. In particular, the invention employs a suite of non-linear mapping algorithms to represent chemical compounds as objects in preferably 2D or 3D Euclidean space.

According to the invention, the distances between objects in that space represent the similarities and/or dissimilarities of the corresponding compounds (relative to selected properties or features of the compounds) computed by some prescribed method. The resulting maps are displayed on a suitable graphics device (such as a graphics terminal, for example), and interactively analyzed to reveal relationships between the data, and to initiate an array of tasks related to these compounds.

2. *Structure of the Invention*

FIG. 1 is a block diagram of a computing environment 102 according to a preferred embodiment of the present invention.

A chemical data visualization and interactive analysis module 104 includes a map generating module 106 and user interface modules 108. The map generating module 106 determines distances between chemical compounds relative to one or more selected properties or features (herein sometimes called evaluation properties or features) of the compounds. The map generating module 106 performs this function by retrieving and analyzing data on chemical compounds and reagents from reagent and compound databases 122. These reagent and compound databases 122 store information on chemical compounds and reagents of interest.

The reagent and compound databases 122 are part of databases 120, which communicate with the chemical data visualization and interactive analysis module 104

via a communication medium 118. The communication medium 118 is preferably any type of data communication means, such as a data bus, a computer network, etc.

The user interface modules 108, which include a map viewer 112 and optionally a structure browser 110, displays a preferably 2D or 3D non-linear map on a suitable graphics device. The non-linear map includes objects that represent the chemical compounds, where the distances between the objects in the non-linear map are those distances determined by the map generating module 106. The user interface modules 108 enable human operators to interactively analyze and process the information in the non-linear map so as to reveal relationships between the data, and to initiate an array of tasks related to the corresponding compounds.

The user interface modules 108 enable users to organize compounds as collections (representing, for example, a combinatorial library). Information pertaining to compound collections are preferably stored in a collection database 124. Information on reagents that are mixed to form compound collections are preferably stored in a library database 126.

Input Device(s) 114 receive input (such as data, commands, queries, etc.) from human operators and forward such input to, for example, the chemical data visualization and interactive analysis module 104 via the communication medium 118. Any well known, suitable input device can be used in the present invention, such as a keyboard, pointing device (mouse, roller ball, track ball, light pen, etc.), touch screen, voice recognition, etc. User input can also be stored and then retrieved, as appropriate, from data/command files.

Output Device(s) 116 output information to human operators. Any well known, suitable output device can be used in the present invention, such as a monitor, a printer, a floppy disk drive or other storage device, a text-to-speech synthesizer, etc.

As described below, the present invention enables the chemical data visualization and interactive analysis module 104 to interact with a number of other modules, including but not limited to one or more map viewers 112, NMR (nuclear magnetic resonance) widget/module 130, structure viewers 110, MS (mass spectrometry) widget/module 134, spreadsheets 136, QSAR (Quantitative Structure-

Activity Relationships) module 138, an experiment planner 140, property prediction programs 142, active site docker 144, etc. These modules communicate with the chemical data visualization and interactive analysis module 104 via the communication medium 118.

3. *Implementation Embodiment of the Invention*

Components shown in the computing environment 102 of FIG. 1 (such as the chemical data visualization and interactive analysis module 104) can be implemented using one or more computers, such as an example computer 202 shown in FIG. 2.

The computer 202 includes one or more processors, such as processor 204. Processor 204 is connected to a communication bus 206. Various software embodiments are described in terms of this example computer system. After reading this description, it will become apparent to a person skilled in the relevant art(s) how to implement the invention using other computer systems and/or computer architectures.

Computer 202 also includes a main memory 208, preferably random access memory (RAM), and can also include one or more secondary storage devices 210. Secondary storage devices 210 can include, for example, a hard disk drive 212 and/or a removable storage drive 214, representing a floppy disk drive, a magnetic tape drive, an optical disk drive, etc. Removable storage drive 214 reads from and/or writes to a removable storage unit 216 in a well known manner. Removable storage unit 216 represents a floppy disk, magnetic tape, optical disk, etc. which is read by and written to by removable storage drive 214. Removable storage unit 216 includes a computer usable storage medium having stored therein computer software and/or data.

In alternative embodiments, the computer 202 can include other similar means for allowing computer programs or other instructions to be loaded into computer 202. Such means can include, for example, a removable storage unit 220 and an interface 218. Examples of such can include a program cartridge and cartridge interface (such as that found in video game devices), a removable memory chip (such as an EPROM,

or PROM) and associated socket, and other removable storage units 220 and interfaces 218 which allow software and data to be transferred from the removable storage unit 220 to computer 202.

The computer 202 can also include a communications interface 222. Communications interface 222 allows software and data to be transferred between computer 202 and external devices. Examples of communications interface 222 include, but are not limited to a modem, a network interface (such as an Ethernet card), a communications port, a PCMCIA slot and card, etc. Software and data transferred via communications interface 222 are in the form of signals which can be electronic, electromagnetic, optical or other signals capable of being received by communications interface 222.

In this document, the term "computer program product" is used to generally refer to media such as removable storage units 216, 220, a hard drive 212 that can be removed from the computer 202, and signals carrying software received by the communications interface 222. These computer program products are means for providing software to the computer 202.

Computer programs (also called computer control logic) are stored in main memory and/or secondary storage devices 210. Computer programs can also be received via communications interface 222. Such computer programs, when executed, enable the computer 202 to perform the features of the present invention as discussed herein. In particular, the computer programs, when executed, enable the processor 204 to perform the features of the present invention. Accordingly, such computer programs represent controllers of the computer 202.

In an embodiment where the invention is implemented using software, the software can be stored in a computer program product and loaded into computer 202 using removable storage drive 214, hard drive 212, and/or communications interface 222. The control logic (software), when executed by the processor 204, causes the processor 204 to perform the functions of the invention as described herein.

In another embodiment, the automated portion of the invention is implemented primarily in hardware using, for example, hardware components such as application

specific integrated circuits (ASICs). Implementation of the hardware state machine so as to perform the functions described herein will be apparent to persons skilled in the relevant art(s).

In yet another embodiment, the invention is implemented using a combination of both hardware and software.

The computer 202 can be any suitable computer, such as a computer system running an operating system supporting a graphical user interface and a windowing environment. A suitable computer system is a Silicon Graphics, Inc. (SGI) workstation/server, a Sun workstation/server, a DEC workstation/server, an IBM workstation/server, an IBM compatible PC, an Apple Macintosh, or any other suitable computer system, such as one using one or more processors from the Intel Pentium family, such as Pentium Pro or Pentium II. Suitable operating systems include, but are not limited to, IRIX, OS/Solaris, Digital Unix, AIX, Microsoft Windows 95/NT, Apple Mac OS, or any other operating system supporting a graphical user interface and a windowing environment. For example, in a preferred embodiment the program may be implemented and run on an Silicon Graphics Octane workstation running the IRIX 6.4 operating system, and using the Motif graphical user interface based on the X Window System.

4. *Overview of Multidimensional Scaling (MDS) and Non-Linear Mapping (NLM)*

According to the present invention, multidimensional scaling (MDS) and non-linear mapping (NLM) techniques are used to generate the non-linear map (i.e., the non-linear map) that includes objects, where the objects represent chemical compounds, and the distances between the objects are indicative of the similarities and dissimilarities between the corresponding compounds. MDS and NLM are described in this section.

MDS and NLM were introduced by Torgerson, *Psychometrika*, 17:401 (1952); Kruskal, *Psychometrika*, 29:115 (1964); and Sammon, *IEEE Trans. Comput.*,

C-18:401 (1969) as a means to generate low-dimensional representations of psychological data. Multidimensional scaling and non-linear mapping are reviewed in Schiffman, Reynolds and Young, *Introduction to Multidimensional Scaling*, Academic Press, New York (1981); Young and Hamer, *Multidimensional Scaling: History, Theory and Applications*, Erlbaum Associates, Inc., Hillsdale, NJ (1987); and Cox and Cox, *Multidimensional Scaling*, Number 59 in *Monographs in Statistics and Applied Probability*, Chapman-Hall (1994). The contents of these publications are incorporated herein by reference in their entireties.

4.1 Procedure Suitable for Relatively Small Data Sets

MDS and NLM (these are generally the same, and are hereafter collectively referred to as MDS) represent a collection of methods for visualizing proximity relations of objects by distances of points in a low-dimensional Euclidean space. Proximity measures are reviewed in Hartigan, *J. Am. Statist. Ass.*, 62:1140 (1967), which is incorporated herein by reference in its entirety. In particular, given a finite set of vectorial or other samples $A = \{a_i, i = 1, \dots, k\}$, a distance function $d_{ij} = d(a_i, a_j)$, with $a_i, a_j \in A$, which measures the similarity and dissimilarity between the i -th and j -th objects in A , and a set of images $X = \{x_i, \dots, x_k; x_i \in \mathbb{R}^m\}$ of A on an m -dimensional display plane (\mathbb{R}^m being an m dimensional vector of real numbers), the objective is to place x_i onto the display plane in such a way that their Euclidean distances $\|x_i - x_j\|$ approximate as closely as possible the corresponding values d_{ij} . This projection, which can only be made approximately, is carried out in an iterative fashion by minimizing an error function which measures the difference between the distance matrices of the original and projected vector sets. Several such error functions have been proposed, most of which are of the least-squares type, including Kruskal's 'stress':

$$S = \sqrt{\frac{\sum_{i < j}^k (d_{ij} - \delta_{ij})^2}{\sum_{i < j}^k d_{ij}^2}} \quad \text{EQ. 1}$$

Sammon's error criterion:

$$E = \frac{\sum_{i < j}^k \frac{(d_{ij} - \delta_{ij})^2}{d_{ij}}}{\sum_{i < j}^k d_{ij}} \quad \text{EQ. 2}$$

and Lingoes' alienation coefficient:

$$K = \sqrt{\frac{\sum_{i < j}^k (d_{ij} \delta_{ij})^2}{\sum_{i < j}^k \delta_{ij}}} \quad \text{EQ. 3}$$

where $\delta_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$ is the Euclidean distance between the images \mathbf{x}_i and \mathbf{x}_j on the display plane. Generally, the solution is found in an iterative fashion by (1) computing or retrieving from a database the distances d_{ij} ; (2) initializing the images \mathbf{x}_i ; (3) computing the distances of the images δ and the value of the error function (e.g. S, E or K in EQ. 1-3 above); (4) computing a new configuration of the images \mathbf{x}_i using a gradient descent procedure, such as Kruskal's linear regression or Guttman's rank-image permutation; and (5) repeating steps 3 and 4 until the error is minimized within some prescribed tolerance.

For example, the Sammon algorithm minimizes EQ. 2 by iteratively updating the coordinates x , using Eq 4:

$$x_{pq}(m+1) = x_{pq}(m) - \lambda \Delta_{pq}(m) \quad \text{EQ. 4}$$

where m is the iteration number, x_{pq} is the q -th coordinate of the p -th image x , λ is the learning rate, and

$$\Delta_{pq}(m) = \frac{\frac{\partial E(m)}{\partial x_{pq}(m)}}{\sqrt{\frac{\partial^2 E(m)}{\partial x_{pq}(m)^2}}} \quad \text{EQ. 5}$$

The partial derivatives in EQ. 5 are given by:

$$\frac{\partial E(m)}{\partial x_{pq}(m)} = -2 \frac{\sum_{j=1, j \neq p}^k \frac{d_{pj} - \delta_{pj}}{d_{pj} \delta_{pj}} (x_{pq} - x_{jq})}{\sum_{i < j} d_{ij}} \quad \text{EQ. 6}$$

$$\frac{\partial^2 E(m)}{\partial x_{pq}(m)^2} = -2 \frac{\sum_{i < j} \frac{1}{d_{ij} \delta_{ij}} \left| (d_{ij} - \delta_{ij}) - \frac{(x_{pq} - x_{jq})^2}{\delta_{ij}} \left(1 + \frac{(d_{ij} - \delta_{ij})}{\delta_{ij}} \right) \right|}{\sum_{i < j} d_{ij}} \quad \text{EQ. 7}$$

The non-linear mapping is obtained by repeated evaluation of EQ. 2, followed by modification of the coordinates using EQ. 4 and 5, until the error is minimized within a prescribed tolerance.

4.2 Procedure Suitable for Large Data Sets

The general refinement paradigm described in Section 4.1 is suitable for relatively small data sets, but has one important limitation that renders it impractical for large data sets. This limitation stems from the fact that the computational effort required to compute the gradients scales to the square of the size of the data set. For relatively large data sets, this quadratic time complexity makes even a partial refinement intractable.

According to the present invention, the following approach is used for large data sets. This approach is to use iterative refinement based on 'instantaneous' errors. As in the approach described in Section 4.1, this approach of Section 4.2 starts with an initial configuration of points generated at random or by some other procedure (as described below in Section 7). This initial configuration is then continuously refined by repeatedly selecting two points i, j , at random, and modifying their coordinates on the non-linear map according to Eq. 8:

$$x_i(t+1) = f(t, x_i(t), x_j(t), d_{ij}) \quad \text{EQ. 8}$$

where t is the current iteration, $x_i(t)$ and $x_j(t)$ are the current coordinates of the i -th and j -th points on the non-linear map, $x_i(t+1)$ are the new coordinates of the i -th point on the non-linear map, and d_{ij} is the true distance between the i -th and j -th points that we attempt to approximate on the non-linear map (see above). $f(\cdot)$ in EQ. 8 above can assume any functional form. Ideally, this function should try to minimize the difference between the actual and target distance between the i -th and j -th points. For example, $f(\cdot)$ may be given by EQ. 9:

$$\mathbf{x}_i(t+1) = f(t, \mathbf{x}_i(t), \mathbf{x}_j(t), d_{ij}) = \mathbf{x}_i(t) + 0.5 \lambda(t) \frac{(d_{ij} - \delta_{ij}(t))}{\delta_{ij}(t)} (\mathbf{x}_j(t) - \mathbf{x}_i(t))$$

EQ. 9

where t is the iteration number, $\delta_{ij} = \|\mathbf{x}_i(t) - \mathbf{x}_j(t)\|$, and $\lambda(t)$ is an adjustable parameter, referred to hereafter as the 'learning rate.'

An analogous equation has been suggested by Kohonen for the training of self-organizing maps (Kohonen, *Self-Organizing Maps*, Springer-Verlag, Berlin (1995)), incorporated herein by reference in its entirety. This process is repeated for a fixed number of cycles, or until some global error criterion is minimized within some prescribed tolerance. A large number of iterations are typically required to achieve statistical accuracy.

The method described above is generally reminiscent of Kohonen's self-organizing principle (Kohonen, *Biological Cybernetics*, 43:59 (1982)) and neural network back-propagation training (Werbos, *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*, PhD Thesis, Harvard University, Cambridge, MA (1974)), and Rumelhart and McClelland, Eds., *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, Vol. 1, MIT Press, Cambridge, MA (1986)), all of which are incorporated herein by reference in their entirety.

The learning rate $\lambda(t)$ in EQ. 9 plays a key role in ensuring convergence. If λ is too small, the coordinate updates are small, and convergence is slow. If, on the other hand, λ is too large, the rate of learning may be accelerated, but the non-linear map may become unstable (i.e. oscillatory). Typically, λ ranges in the interval $[0, 1]$ and may be fixed, or it may decrease monotonically during the refinement process. Moreover, λ may also be a function of i, j and/or d_{ij} , and can be used to apply different weights to certain objects, distances and/or distance pairs. For example, λ may be computed by EQ. 10:

$$\lambda(t) = (\lambda_{\max} + t \frac{\lambda_{\min} - \lambda_{\max}}{T}) \frac{1}{1 + a d_{ij}}$$

EQ. 10

or EQ. 11:

$$\lambda(t) = (\lambda_{\max} + t \frac{\lambda_{\min} - \lambda_{\max}}{T}) e^{-a d_{ij}}$$

EQ. 11

where λ_{\max} and λ_{\min} are the (unweighted) starting and ending learning rates such that $\lambda_{\max}, \lambda_{\min} \in [0,1]$, T is the total number of refinement steps (iterations), t is the current iteration number, and a is a constant scaling factor. EQ. 10 and 11 have the effect of decreasing the correction at large separations, thus creating a non-linear map which preserves short-range interactions more faithfully than long-range ones. Weighting is discussed in greater detail below. Because of the general resemblance of the training process described above to Kohonen's self-organizing principle, these maps shall sometimes be herein called 'Self-Organizing Non-Linear Maps.'

One of the main advantages of this approach is that it makes partial refinements possible. It is often sufficient that the pair-wise dissimilarities are represented only approximately to reveal the general structure and topology of the data. Unlike traditional MDS, this approach allows very fine control of the refinement process. Moreover, as the non-linear map self-organizes, the pair-wise refinements become cooperative, which partially alleviates the quadratic nature of the problem.

The general usefulness of multi-dimensional scaling stems from the fact that data in \mathbb{R}^d are almost never d -dimensional. Although scaling becomes more problematic as

the *true* dimensionality of the space increases, the presence of structure in the data is very frequently reflected on the resulting map. Of course, one can easily conceive of situations where MDS is not effective, particularly when the data is random and truly hyper-dimensional. Fortunately, these situations rarely arise in practice, as some form of structure is always present in the data, particularly data related to molecular structure and function.

The embedding procedure described above does not guarantee convergence to the global minimum (i.e., the most faithful embedding in a least-squares sense). If so desired, the refinement process may be repeated a number of times from different starting configurations and/or random number seeds. It should also be pointed out that the absolute coordinates in the non-linear map carry no physical significance. What is important are the relative distances between points, and the general structure and topology of the data (presence, density and separation of clusters, etc.).

The method described above is ideally suited for both metric and non-metric scaling. The latter is particularly useful when the (dis)similarity measure is not a true metric, i.e. it does not obey the distance postulates and, in particular, the triangle inequality (such as the Tanimoto coefficient, for example). Although an 'exact' projection is only possible when the distance matrix is positive definite, meaningful projections can still be obtained even when this criterion is not satisfied. As mentioned above, the overall quality of the projection is determined by a sum-of-squares error function such as those shown in EQ. 1-3.

5. *Evaluation Properties (Features) and Distance Measures*

As mentioned above, the distances d_{ij} between chemical compounds are computed according to some prescribed measure of molecular 'similarity'. This similarity can be based on any combination of properties or features of the compounds. For example, the similarity measure may be based on structural similarity, chemical similarity, physical similarity, biological similarity, and/or some other type of similarity measure which can be derived from the structure or identity of

the compounds. Under the system of the present invention, any similarity measure can be used to construct the non-linear map. The properties or features that are being used to evaluate similarity or dissimilarity among compounds are sometimes herein collectively called "evaluation properties."

5.1 *Evaluation Properties Having Continuous or Discrete Real Values*

As noted above, in a preferred embodiment of the present invention, the similarity measure may be derived from a list of physical, chemical and/or biological properties (i.e., evaluation properties) associated with a set of compounds. Under this formalism, the compounds are represented as vectors in multi-variate property space, and their similarity may be computed by some geometrical distance measure.

In a preferred embodiment, the property space is defined using one or more molecular features (descriptors). Such molecular features may include topological indices, physicochemical properties, electrostatic field parameters, volume and surface parameters, etc. For example, these features may include, but are not limited to, molecular volume and surface areas, dipole moments, octanol-water partition coefficients, molar refractivities, heats of formation, total energies, ionization potentials, molecular connectivity indices, 2D and 3D auto-correlation vectors, 3D structural and/or pharmacophoric parameters, electronic fields, etc. However, it should be understood that the present invention is not limited to this embodiment. For example, molecular features may include the observed biological activities of a set of compounds against an array of biological targets such as enzymes or receptors (also known as affinity fingerprints). In fact, any vectorial representation of chemical data can be used in the present invention.

5.2 Distance Measure Where Values of Evaluation Properties Are Continuous or Discrete Real Numbers

A “distance measure” is some algorithm or technique used to determine the difference between compounds based on the selected evaluation properties. The particular distance measure that is used in any given situation depends, at least in part, on the set of values that the evaluation properties can take.

For example, where the evaluation properties can take real numbers as values, then a suitable distance measure is the Minkowski metric, shown in EQ. 12:

$$d_{ij} = d(x_i, x_j) = \left(\sum_k |x_{ik} - x_{jk}|^r \right)^{\frac{1}{r}} \quad \text{EQ. 12}$$

where k is used to index the elements of the property vector, and $r \in [1, \infty)$. For $r = 1.0$, EQ. 12 is the city-block or Manhattan metric. For $r = 2.0$, EQ. 12 is the ordinary Euclidean metric. For $r = \infty$, EQ. 12 is the maximum of the absolute coordinate distances, also referred to as the ‘dominance’ metric, the ‘sup’ metric, or the ‘ultrametric’ distance. For any value of $r \in [1, \infty)$, it can be shown that the Minkowski metric is a true metric, i.e. it obeys the distance postulates and, in particular, the triangle inequality.

5.3 Evaluation Properties Having Binary Values

Alternatively, the evaluation properties of the compounds may be represented in a binary form (i.e., either a compound has or does not have an evaluation property), where each bit is used to indicate the presence or absence (or potential presence or absence) of some molecular feature or characteristic. For example, compounds may be encoded using substructure keys where each bit is used to denote the presence or absence of a specific structural feature or pattern in the target molecule. Such features

include, but are not limited to, the presence, absence or minimum number of occurrences of a particular element (e.g. the presence of at least 1, 2 or 3 nitrogen atoms), unusual or important electronic configurations and atom types (e.g. doubly-bonded nitrogen or aromatic carbon), common functional groups such as alcohols, amines *etc.*, certain primitive and composite rings, a pair or triplet of pharmacophoric groups at a particular separation in 3-dimensional space, and 'disjunctions' of unusual features that are rare enough not to worth an individual bit, yet extremely important when they do occur (typically, these unusual features are assigned a common bit that is set if any one of the patterns is present in the target molecule).

Alternatively, the evaluation properties of compounds may be encoded in the form of binary fingerprints, which do not depend on a predefined fragment or feature dictionary to perform the bit assignment. Instead, every pattern in the molecule up to a predefined limit is systematically enumerated, and serves as input to a hashing algorithm that turns 'on' a small number of bits at pseudo-random positions in the bitmap. Although it is conceivable that two different molecules may have exactly the same fingerprint, the probability of this happening is extremely small for all but the simplest cases. Experience suggests that these fingerprints contain sufficient information about the molecular structures to permit meaningful similarity comparisons.

5.4 Distance Measures Where Values of Evaluation Properties Are Binary

A number of similarity (distance) measures can be used with binary descriptors (i.e., where evaluation properties are binary or binary fingerprints). The most frequently used ones are the normalized Hamming distance:

$$H = \frac{|XOR(x,y)|}{N} \quad \text{EQ. 13}$$

which measures the number of bits that are different between x and y, the Tanimoto or Jaccard coefficient:

$$T = \frac{|AND(x,y)|}{|OR(x,y)|} \quad \text{EQ. 14}$$

which is a measure of the number of substructures shared by two molecules relative to the ones they *could* have in common, and the Dice coefficient:

$$D = \frac{2|AND(x,y)|}{|x| + |y|} \quad \text{EQ. 15}$$

In the equations listed above, AND(x, y) is the intersection of binary sets x and y (bits that are 'on' in both sets), IOR(x, y) is the union or 'inclusive or' of x and y (bits that are 'on' in either x or y), XOR is the 'exclusive or' of x and y (bits that are 'on' in either x or y, but not both), |x| is the number of bits that are 'on' in x, and N is the length of the binary sets measured in bits (a constant).

Another popular metric is the Euclidean distance which, in the case of binary sets, can be recast in the form:

$$E = \sqrt{N - |XOR(x, NOT(y))|} \quad \text{EQ. 16}$$

where NOT(y) denotes the binary complement of y. The expression |XOR(x, NOT(y))| represents the number of bits that are identical in x and y (either 1's or 0's). The Euclidean distance is a good measure of similarity when the binary sets are relatively rich, and is mostly used in situations in which similarity is measured in a relative sense.

In the examples described above, the distance between two compounds is determined using a binary or multivariate representation. However, the system of the

present invention is not limited to this embodiment. For example, the similarity between two compounds may be determined by comparing the shapes of the molecules using a suitable 3-dimensional alignment method, or it may be inferred by a similarity model defined according to a prescribed procedure. For example, one such similarity model may be a neural network trained to predict a similarity coefficient given a suitably encoded pair of compounds. Such a neural network may be trained using a training set of structure pairs and a known similarity coefficient for each such pair, as determined by user input, for example.

6. *Scaling of Evaluation Properties*

Referring back to EQ. 12, according to the present invention, the features (i.e., evaluation properties) may be scaled differently to reflect their relative importance in assessing the proximity between two compounds. For example, suppose the user has selected two evaluation properties, Property A and Property B. If Property A has a weight of 2, and Property B has a weight of 10, then Property B will have five times the impact on the distance calculation than Property A.

According to this embodiment of the invention, EQ. 12 may be replaced by EQ. 17:

$$d_y = d(x_i, x_j) = \left(\sum_k (w_k |x_{ik} - x_{jk}|^r) \right)^{\frac{1}{r}} \quad \text{EQ. 17}$$

where w_k is the weight of the k-th property. An example of such a weighting factor is a normalization coefficient. However, other weighting schemes may also be used.

According to the present invention, the scaling (weights) need not be uniform throughout the entire map, i.e. the resulting map need not be isomorphic. Hereafter, maps derived from uniform weights shall be referred to as globally weighted (isomorphic), whereas maps derived from non-uniform weights shall be referred to as

locally weighted (non-isomorphic). On locally-weighted maps, the distances on the non-linear map reflect a local measure of similarity. That is, what determines similarity in one domain of the non-linear map is not necessarily the same with what determines similarity on another domain of the non-linear map. For example, locally-weighted maps may be used to reflect similarities derived from a locally-weighted case-based learning algorithm. Locally-weighted learning uses locally weighted training to average, interpolate between, extrapolate from, or otherwise combine training data. Most learning methods (also referred to as modeling or prediction methods) construct a single model to fit all the training data. Local models, on the other hand, attempt to fit the training data in a local region around the location of the query. Examples of local models include nearest neighbors, weighted average, and locally weighted regression. Locally-weighted learning is reviewed in Vapnik, in *Advances in Neural Information Processing Systems*, 4:831, Morgan-Kaufman, San Mateo, CA (1982); Bottou and Vapnik, *Neural Computation*, 4(6):888 (1992); and Vapnik and Bottou, *Neural Computation*, 5(6):893 (1993), all of which are incorporated herein by reference in their entireties.

According to the present invention, it is also possible to construct a non-linear map from a distance matrix which is not strictly symmetric, i.e. a distance matrix where $d_{ij} \neq d_{ji}$. A potential use of this approach is in situations where the distance function is defined locally, e.g. in a locally weighted model using a point-based local distance function. In this embodiment, each training case has associated with it a distance function and the values of the corresponding parameters. Preferably, to construct a non-linear map which reflects these local distance relationships, the distance between two points is evaluated twice, using the local distance functions of the respective points. The resulting distances are averaged, and are used as input in the non-linear mapping algorithm described above. If the point-based local distance functions vary in some continuous or semi-continuous fashion throughout the feature space, this approach could potentially lead to a meaningful projection.

7. Improvements to Map Generation Process

This section describes improvements to the chemical visualization map generation process described above. Each of the enhancements described below is under the control of the user. That is, the user can elect to perform or not perform each of the enhancements discussed below. Alternatively, the invention can be defined so that the below enhancements are automatically performed, unless specifically overridden by the user (or in some embodiments, the user may not have the option of overriding one or more of the below enhancements).

7.1 Pre-Ordering

In many cases, the approach described above for generating the non-linear map may be accelerated by pre-ordering the data using a suitable statistical method. For example, if the data is available in vectorial or binary form, the initial configuration of the points on the non-linear map may be computed using Principal Component Analysis. In a preferred embodiment, the initial configuration may be constructed from the first 3 principal components of the feature matrix (i.e. the 3 latent variables which account for most of the variance in the data). In practice, this technique can have profound effects in the speed of refinement. Indeed, if a random initial configuration is used, a significant portion of the training time is spent establishing the general structure and topology of the non-linear map, which is typically characterized by large rearrangements. If, on the other hand, the input configuration is partially ordered, the error criterion can be reduced relatively rapidly to an acceptable level.

7.2 Localized Refinement

If the data is highly clustered, by virtue of the sampling process low-density areas may be refined less effectively than high-density areas. In a preferred embodiment, this tendency may be partially compensated by a modification to the original algorithm which increases the sampling probability in low-density areas. In

one embodiment, the center of mass of the non-linear map is identified, and concentric shells centered at that point are constructed. A series of regular refinement iterations are then carried out, each time selecting points from within or between these shells. This process is repeated for a prescribed number of cycles. This phase is then followed by a phase of regular refinement using global sampling, and the process is repeated.

As mentioned above, the basic algorithm does not distinguish short- from long-range distances. EQ. 10 and 11 describe a method to ensure that short-range distances are preserved more faithfully than long-range ones through the use of weighting. An alternative (and complementary) approach is to ensure that points at close separation are sampled more extensively than points at long separation. A preferred embodiment is to use an alternating sequence of global and local refinement cycles, similar to the one described above. In this embodiment, a phase of global refinement is initially carried out. At the end of this phase, the resulting non-linear map is partitioned into a regular grid, and the points (objects) in each cell are subjected to a phase of local refinement (i.e. only points from within the same cell are compared and refined). Preferably, the number of sampling steps in each cell should be proportional to the number of points contained in that cell. This process is highly parallelizable. This local refinement phase is then followed by another global refinement phase, and the process is repeated for a prescribed number of cycles, or until the embedding error is minimized within a prescribed tolerance. Alternatively, the grid method may be replaced by another suitable method for identifying proximal points, such as a k-d tree, for example.

7.3 Incremental Refinement

The approach and techniques described herein may be used for incremental refinement of a map. That is, starting from an organized non-linear map of a set of objects or points (compounds), a new set of objects (compounds) may be added without modification of the original map. Strictly speaking, this is statistically

acceptable if the new set of objects is significantly smaller than the original set. In a preferred embodiment, the new set of objects may be 'diffused' into the existing map, using a modification of the algorithm described above. In particular, EQ. 8 and 9 can be used to update only the new objects. In addition, the sampling procedure ensures that the selected pairs contain at least one object from the incoming set. That is, two objects are selected at random so that at least one of these objects belongs to the incoming set.

8. *Operation of the Present Invention*

The operation of the present invention with regard to visualizing and interactively processing chemical compounds in a non-linear map shall now be described with reference to a flowchart 302 shown in FIG. 3. Unless otherwise specified, interaction with users described below is achieved by operation of the user interface modules 108 (FIG. 1).

In step 304, the user selects one or more compounds to map in a new non-linear map. The user may select compounds to map by retrieving a list of compounds from a file, by manually typing in a list of compounds, and/or by using a graphical user interface (GUI) such as the structure browser shown in FIG. 5 (described below). The invention envisions other means for enabling the user to specify compounds to display in a non-linear map. For example, the user can also select compounds from an already existing compound visualization non-linear map (in one embodiment, the user drags and drops the compounds from the old compound visualization non-linear map to the new compound visualization non-linear map -- drag and drop operations according to the present invention are described below).

In step 306, the user selects a method to be used for evaluating the molecular similarity or dissimilarity between the compounds selected in step 304. In an embodiment, the similarity/dissimilarity between the compounds selected in step 304 is determined (in step 308) based on a prescribed set of evaluation properties. As described above, evaluation properties can be any properties related to the structure,

function, or identity of the compounds selected in step 304. Evaluation properties include, but are not limited to, structural properties, functional properties, chemical properties, physical properties, biological properties, etc., of the compounds selected in step 304.

5 In an embodiment of the present invention, the selected evaluation properties may be scaled differently to reflect their relative importance in assessing the proximity (i.e., similarity or dissimilarity) between two compounds. Accordingly, also in step 306, the user selects a scale factor for each of the selected evaluation. Note that such selection of scale factors is optional. The user need not select a scale factor for each selected evaluation property. If the user does not select a scale factor for a given evaluation property, then that evaluation property is given a default scale factor, such as unity.

Alternatively in step 306, the user can elect to retrieve similarity/dissimilarity values pertaining to the compounds selected in step 304 from a source, such as a database. These similarity/dissimilarity values in the database were previously generated. In another embodiment, the user in step 306 can elect to determine similarity/dissimilarity values using any well-known technique or procedure.

10 In step 308, the map generating module 106 generates a new non-linear map. This new non-linear map includes a point for each of the compounds selected in step 304. Also, in this new non-linear map, the distance between any two points is representative of their similarity/dissimilarity. The manner in which the map generating module 106 generates the new non-linear map shall now be further described with reference to a flowchart 402 in FIG. 4.

20 In step 404, coordinates on the new non-linear map of points corresponding to the compounds selected in step 304 are initialized.

In step 406, two of the compounds i, j selected in step 304 are selected for processing.

25 In step 408, similarity/dissimilarity d_{ij} between compounds i, j is determined based on the method selected by the user in step 306.

In step 410, based on the similarity/dissimilarity d_{ij} determined in step 408, coordinates of points corresponding to compounds i, j on the non-linear map are obtained.

In step 412, training/learning parameters are updated.

5 In step 414, a decision is made as to terminate or not terminate. If a decision is made to not terminate at this point, then control returns to step 406. Otherwise, step 416 is performed.

In step 416, the non-linear map is output (i.e., generation of the non-linear map is complete).

Details regarding the steps of flowchart 402 are discussed above.

Referring again to FIG. 3, in step 312 the map viewer 112 displays the new non-linear map on an output device 116 (such as a computer graphics monitor). Examples of non-linear maps being displayed by the map viewer 112 are shown in FIGS. 6 and 7 (described below).

In step 314, the user interface modules 108 enable operators to interactively analyze and process the compounds represented in the displayed non-linear map. These user interface functions of the present invention are described below.

20 The present invention enables users to modify existing compound visualization non-linear maps (as used herein, the term "compound visualization non-linear map" refers to a rendered non-linear map). For example, users can add additional compounds to the map, remove compounds from the map, highlight compounds on the map, etc. In such cases, pertinent functional steps of flowchart 302 are repeated. For example, steps 304 (selecting compounds to map), 310 (generating the non-linear map), and 312 (displaying the map) are repeated when the user opts to
25 add new compounds to an existing map. However, according to an embodiment of the invention, the map is incrementally refined and displayed in steps 310 and 312 when adding compounds to an existing compound visualization non-linear map (this incremental refinement is described above).

9. User Interface of the Present Invention

The user interface features of the present invention are described in this section. Various user interface modules and features are described below. Also, various functional/control threads (in the present context, a functional/control thread is a series of actions performed under the control of a user) employing these user interface modules and features are described below. It will be appreciated by persons skilled in the relevant art(s) that the user interface of the present invention is very flexible, varied, and diverse. An operator can employ the user interface of the present invention to perform a wide range of activities with respect to visualizing and interactively analyzing chemical compounds. Accordingly, it should be understood that the functional/control threads described herein are provided for illustrative purposes only. The invention is not limited to these functional/control threads.

Preferably, the invention provides the following capabilities, features, and functions: displaying 2D and/or 3D chemical structures and/or chemical names; displaying compound collections and/or libraries; displaying components of structures (i.e. building blocks) of combinatorial libraries; visualization of compound collections and/or libraries as 2D and/or 3D maps of colored objects.

Also, the present invention allows the following: (1) browsing compound collections and/or libraries; (2) selection of individual compounds, collections of compounds and/or libraries of compounds; (3) selection of compounds generated in a combinatorial fashion via selection of their respective building blocks; (4) mapping, visualization, and/or linking of compounds onto and/or from 2D and/or 3D maps; (5) manipulation of the 2D and/or 3D maps such as rotation, resizing, translation, etc.; (6) manipulation of objects on the 2D and/or 3D maps such as changing the appearance of objects (visibility, size, shape, color, etc.), changing position of objects on the map, and/or changing relationships between objects on the map; (7) interactive exploring of the 2D and/or 3D maps such as querying chemical structure, querying distance, selection of individual objects and/or areas of a map, etc.

Additional user interface features, functions, and capabilities of the present invention will be apparent to persons skilled in the relevant art(s) based on the discussion contained herein.

As shown in FIG. 1, the invention includes a structure browser 110 and a map viewer 112. At any given time, each of these can have multiple instances depending on the program use.

9.1 Structure Browser

FIG. 5 illustrates a structure browser window 502 generated by the structure browser 110. The structure browser window 502 includes a frame 504, a menu pane 506, and a group of labeled tabbed pages 508. Each tabbed page holds a molecular spreadsheet or a group of labeled tabbed pages.

Each tab is associated with a compound collection (tabs 510) or a library, such as a combinatorial library (tabs 512). Selecting a collection tab 510 brings up a table of corresponding chemical structures. Selecting a library tab 512 brings up a group of tabbed pages corresponding to the sets of building blocks used to generate the library. Each of the library's tabbed pages works the same way as a compound collection tabbed page. In the example shown in FIG. 5, the tab 510 called "DDL0" is selected. DDL0 has three building block tabs 512, called "Cores," "Acids," and "Amines." The "Acids" collection tab is currently selected, so that a table 522 of the structures of the compounds in the "Acids" collection is shown.

The browser window 502 includes a table 522, a slider 514, an input field 516, and two buttons: "Prev Page" 518 and "Next Page" 520. The slider 514, the input field 516, and the buttons 518, 520 facilitate browsing the content of the Acids table 522. If we consider the content of the table 522 as a contiguous ordered *list* of chemical structures (compounds or building blocks), that shown in the browser window 502 can be considered as a *window* positioned over the *list*. At any given moment this *window* displays part of the *list* depending on its position and the displayed part is equal to the size of the window, i.e., the number of cells in the table. Initially that *window* displays the top of the *list*. Moving the slider 514 changes the position of the *window* over the *list*. Entering a value into the input field 516 specifies the position of the *window* over the *list*. Pushing the "Next Page" button 520 moves

the *window* one window size down the *list*, pushing the “Prev Page” button 518 moves the *window* one window size up the *list*.

The user can select compounds shown in the table 522 for various actions. For example, compounds can be selected using the browser window 502 as input for the generation of a new compound visualization non-linear map, or as input for adding compounds to an existing compound visualization non-linear map. Clicking with a left mouse button over a table cell selects or deselects the corresponding compound structure (toggling). Toggling on/off also changes the color of the cell, to indicate which cells have been selected. Selected structures are displayed on a first background color, and non-selected structures are displayed on a second background color. In the example of FIG. 5, certain cells 523 in table 522 have been selected.

The menu pane 506 contains menus: File, Edit, Selection, Map, and/or other menus. The File menu facilitates file open/save, print, and exit operations. Edit menu contains commands for editing content of the table 522. The Selection menu provides options to select/deselect (clear) a current compound collection, a collection of building blocks of a combinatorial library, and/or all compounds. The Map menu includes commands for creating a map viewer and for displaying a selection of compounds in that map viewer. The latter option brings up a dialog window (FIG. 8), which allows the user to specify shape, color, and/or size of the selected objects, which will be used to represent the selected compounds on the map.

9.2 Map Viewer

A map viewer window 600 generated by the map viewer 112 is shown in FIG. 6. (also see FIGS. 6-10 and 13). A compound visualization non-linear map is displayed in a render area 614 of the map view window 600.

In a preferred embodiment, the map viewer 112 is based on Open Inventor, a C++ library of objects and methods for interactive 3D graphics, publicly available from Silicon Graphics Inc. Open Inventor relies on OpenGL for fast and flexible rendering of 3D objects. Alternatively, the map viewer 112 can be based on a publicly

available VRML viewer. Alternatively, any other software and/or hardware product allowing rendering of 3D objects/scenes can be used.

In a preferred embodiment, 3D compound visualization maps of chemical compounds are implemented as Open Inventor 3D scene databases. Each map is build
5 as an ordered collection of nodes referred to as a scene graph. Each scene graph includes, but is not limited to, nodes representing cameras (points of view), light sources, 3D shapes, objects surface materials, and geometric transformations. Each chemical compound displayed on a map is associated with a 3D shape node, a material node and a geometric transformation node.

Geometric transformation node reflects compound coordinates in the map. 3D shape node and material node determine shape, size and color of the visual object associated with the compound. Combinations of a particular shape, size and color are used to display compounds grouped by a certain criteria, thus allowing easy visual differentiation of different groups/sets of compounds. 3D shapes of the visual objects
10 in the map include, but not limited to, point, cube, sphere, and cone. Color of a visual object in the map can be set to any combination of three basic colors: red, green and blue. Besides the color, material node can specify transparency and shininess of a visual object's surface.

In an embodiment, an object's display properties (color, intensity of color, transparent, degree of transparency, shininess, degree of shininess, etc.) can represent
20 physical, chemical, biological, and/or other properties of the corresponding compound, such as the cost of the compound, difficulty of synthesizing the compound, whether the compound is available in a compound repository, etc. For example, the larger the molecular weight of an object, the larger the size of the corresponding object in the display map.

Each object or point displayed in the compound visualization non-linear map represents a chemical compound. Objects in the compound visualization non-linear map can be grouped into sets.

By default, every time a set of compounds is mapped into a compound
25 visualization non-linear map, a new set of graphical objects is created and added to the

compound visualization non-linear map. All objects in a particular set can share the same attributes: shape, color, and size, thus providing an easy visual identification of the objects belonging to the same set or to different sets.

5 A compound can be a member of several sets. In an embodiment, for a given compound, a different object is displayed in the compound visualization non-linear map for each set of which the compound is a member. In this case the objects in the compound visualization non-linear map that represent the compound as a member of each of the sets may overlap and only the biggest object may be visible. In this case, a toggle sets feature (described below) may be used to reveal multiple set membership.

10 The map viewer window 600 includes a frame 602, a menu pane 604, and a viewer module preferably implemented as an Open Inventor component (examiner viewer). The viewer module incorporates the following elements: (1) a render area 614 in which the compound visualization non-linear map is being displayed; (2) combinations of thumbwheels 608, 610, 612, sliders, and/or viewer functions icons/buttons 620, 622, 624, 626, 628, 630, 632; and (3) pop-up menus and dialogs 15 616, 702, 902 which provide access to all viewers functions, features and/or properties.

20 The thumbwheels 608, 610 rotate the compound visualization non-linear map around a reference point of interest. Thumbwheel 610 rotates in the y direction, and thumbwheel 608 rotates in the x direction. The origin of rotation (i.e., the camera position) is by default the geometric center of the compound visualization map 614 (render area), but can be placed anywhere in the compound visualization non-linear map. The compound visualization non-linear map can also be panned in the screen plane, as well as dollied in and out (forward/backward movement) via thumbwheel 25 612.

The map view window 600 has several different modes or states, e.g. view, pick, panning, dolly, seek, and/or other. Each mode defines a different mouse cursor and how mouse events are interpreted.

In the view mode, mouse motions are translated into rotations of the virtual trackball and corresponding rotations of the compound visualization non-linear map. The view mode is the default mode.

5 In the panning mode, the compound visualization non-linear map is translated in the screen plane following the mouse movements.

In the dolly mode, a scene is moved in and out of screen according to the vertical motions of the mouse.

Seek mode allows the user to change the point of rotation (reference point) of a scene by attaching it to an object displayed in the compound visualization non-linear map.

10
15
20
25
30
35
40
45
50
55
60
65
70
75
80
85
90
95
100
105
110
115
120
125
130
135
140
145
150
155
160
165
170
175
180
185
190
195
200
205
210
215
220
225
230
235
240
245
250
255
260
265
270
275
280
285
290
295
300
305
310
315
320
325
330
335
340
345
350
355
360
365
370
375
380
385
390
395
400
405
410
415
420
425
430
435
440
445
450
455
460
465
470
475
480
485
490
495
500
505
510
515
520
525
530
535
540
545
550
555
560
565
570
575
580
585
590
595
600
605
610
615
620
625
630
635
640
645
650
655
660
665
670
675
680
685
690
695
700
705
710
715
720
725
730
735
740
745
750
755
760
765
770
775
780
785
790
795
800
805
810
815
820
825
830
835
840
845
850
855
860
865
870
875
880
885
890
895
900
905
910
915
920
925
930
935
940
945
950
955
960
965
970
975
980
985
990
1000
1005
1010
1015
1020
1025
1030
1035
1040
1045
1050
1055
1060
1065
1070
1075
1080
1085
1090
1095
1100
1105
1110
1115
1120
1125
1130
1135
1140
1145
1150
1155
1160
1165
1170
1175
1180
1185
1190
1195
1200
1205
1210
1215
1220
1225
1230
1235
1240
1245
1250
1255
1260
1265
1270
1275
1280
1285
1290
1295
1300
1305
1310
1315
1320
1325
1330
1335
1340
1345
1350
1355
1360
1365
1370
1375
1380
1385
1390
1395
1400
1405
1410
1415
1420
1425
1430
1435
1440
1445
1450
1455
1460
1465
1470
1475
1480
1485
1490
1495
1500
1505
1510
1515
1520
1525
1530
1535
1540
1545
1550
1555
1560
1565
1570
1575
1580
1585
1590
1595
1600
1605
1610
1615
1620
1625
1630
1635
1640
1645
1650
1655
1660
1665
1670
1675
1680
1685
1690
1695
1700
1705
1710
1715
1720
1725
1730
1735
1740
1745
1750
1755
1760
1765
1770
1775
1780
1785
1790
1795
1800
1805
1810
1815
1820
1825
1830
1835
1840
1845
1850
1855
1860
1865
1870
1875
1880
1885
1890
1895
1900
1905
1910
1915
1920
1925
1930
1935
1940
1945
1950
1955
1960
1965
1970
1975
1980
1985
1990
2000
2005
2010
2015
2020
2025
2030
2035
2040
2045
2050
2055
2060
2065
2070
2075
2080
2085
2090
2095
2100
2105
2110
2115
2120
2125
2130
2135
2140
2145
2150
2155
2160
2165
2170
2175
2180
2185
2190
2195
2200
2205
2210
2215
2220
2225
2230
2235
2240
2245
2250
2255
2260
2265
2270
2275
2280
2285
2290
2295
2300
2305
2310
2315
2320
2325
2330
2335
2340
2345
2350
2355
2360
2365
2370
2375
2380
2385
2390
2395
2400
2405
2410
2415
2420
2425
2430
2435
2440
2445
2450
2455
2460
2465
2470
2475
2480
2485
2490
2495
2500
2505
2510
2515
2520
2525
2530
2535
2540
2545
2550
2555
2560
2565
2570
2575
2580
2585
2590
2595
2600
2605
2610
2615
2620
2625
2630
2635
2640
2645
2650
2655
2660
2665
2670
2675
2680
2685
2690
2695
2700
2705
2710
2715
2720
2725
2730
2735
2740
2745
2750
2755
2760
2765
2770
2775
2780
2785
2790
2795
2800
2805
2810
2815
2820
2825
2830
2835
2840
2845
2850
2855
2860
2865
2870
2875
2880
2885
2890
2895
2900
2905
2910
2915
2920
2925
2930
2935
2940
2945
2950
2955
2960
2965
2970
2975
2980
2985
2990
2995
3000
3005
3010
3015
3020
3025
3030
3035
3040
3045
3050
3055
3060
3065
3070
3075
3080
3085
3090
3095
3100
3105
3110
3115
3120
3125
3130
3135
3140
3145
3150
3155
3160
3165
3170
3175
3180
3185
3190
3195
3200
3205
3210
3215
3220
3225
3230
3235
3240
3245
3250
3255
3260
3265
3270
3275
3280
3285
3290
3295
3300
3305
3310
3315
3320
3325
3330
3335
3340
3345
3350
3355
3360
3365
3370
3375
3380
3385
3390
3395
3400
3405
3410
3415
3420
3425
3430
3435
3440
3445
3450
3455
3460
3465
3470
3475
3480
3485
3490
3495
3500
3505
3510
3515
3520
3525
3530
3535
3540
3545
3550
3555
3560
3565
3570
3575
3580
3585
3590
3595
3600
3605
3610
3615
3620
3625
3630
3635
3640
3645
3650
3655
3660
3665
3670
3675
3680
3685
3690
3695
3700
3705
3710
3715
3720
3725
3730
3735
3740
3745
3750
3755
3760
3765
3770
3775
3780
3785
3790
3795
3800
3805
3810
3815
3820
3825
3830
3835
3840
3845
3850
3855
3860
3865
3870
3875
3880
3885
3890
3895
3900
3905
3910
3915
3920
3925
3930
3935
3940
3945
3950
3955
3960
3965
3970
3975
3980
3985
3990
3995
4000
4005
4010
4015
4020
4025
4030
4035
4040
4045
4050
4055
4060
4065
4070
4075
4080
4085
4090
4095
4100
4105
4110
4115
4120
4125
4130
4135
4140
4145
4150
4155
4160
4165
4170
4175
4180
4185
4190
4195
4200
4205
4210
4215
4220
4225
4230
4235
4240
4245
4250
4255
4260
4265
4270
4275
4280
4285
4290
4295
4300
4305
4310
4315
4320
4325
4330
4335
4340
4345
4350
4355
4360
4365
4370
4375
4380
4385
4390
4395
4400
4405
4410
4415
4420
4425
4430
4435
4440
4445
4450
4455
4460
4465
4470
4475
4480
4485
4490
4495
4500
4505
4510
4515
4520
4525
4530
4535
4540
4545
4550
4555
4560
4565
4570
4575
4580
4585
4590
4595
4600
4605
4610
4615
4620
4625
4630
4635
4640
4645
4650
4655
4660
4665
4670
4675
4680
4685
4690
4695
4700
4705
4710
4715
4720
4725
4730
4735
4740
4745
4750
4755
4760
4765
4770
4775
4780
4785
4790
4795
4800
4805
4810
4815
4820
4825
4830
4835
4840
4845
4850
4855
4860
4865
4870
4875
4880
4885
4890
4895
4900
4905
4910
4915
4920
4925
4930
4935
4940
4945
4950
4955
4960
4965
4970
4975
4980
4985
4990
4995
5000
5005
5010
5015
5020
5025
5030
5035
5040
5045
5050
5055
5060
5065
5070
5075
5080
5085
5090
5095
5100
5105
5110
5115
5120
5125
5130
5135
5140
5145
5150
5155
5160
5165
5170
5175
5180
5185
5190
5195
5200
5205
5210
5215
5220
5225
5230
5235
5240
5245
5250
5255
5260
5265
5270
5275
5280
5285
5290
5295
5300
5305
5310
5315
5320
5325
5330
5335
5340
5345
5350
5355
5360
5365
5370
5375
5380
5385
5390
5395
5400
5405
5410
5415
5420
5425
5430
5435
5440
5445
5450
5455
5460
5465
5470
5475
5480
5485
5490
5495
5500
5505
5510
5515
5520
5525
5530
5535
5540
5545
5550
5555
5560
5565
5570
5575
5580
5585
5590
5595
5600
5605
5610
5615
5620
5625
5630
5635
5640
5645
5650
5655
5660
5665
5670
5675
5680
5685
5690
5695
5700
5705
5710
5715
5720
5725
5730
5735
5740
5745
5750
5755
5760
5765
5770
5775
5780
5785
5790
5795
5800
5805
5810
5815
5820
5825
5830
5835
5840
5845
5850
5855
5860
5865
5870
5875
5880
5885
5890
5895
5900
5905
5910
5915
5920
5925
5930
5935
5940
5945
5950
5955
5960
5965
5970
5975
5980
5985
5990
5995
6000
6005
6010
6015
6020
6025
6030
6035
6040
6045
6050
6055
6060
6065
6070
6075
6080
6085
6090
6095
6100
6105
6110
6115
6120
6125
6130
6135
6140
6145
6150
6155
6160
6165
6170
6175
6180
6185
6190
6195
6200
6205
6210
6215
6220
6225
6230
6235
6240
6245
6250
6255
6260
6265
6270
6275
6280
6285
6290
6295
6300
6305
6310
6315
6320
6325
6330
6335
6340
6345
6350
6355
6360
6365
6370
6375
6380
6385
6390
6395
6400
6405
6410
6415
6420
6425
6430
6435
6440
6445
6450
6455
6460
6465
6470
6475
6480
6485
6490
6495
6500
6505
6510
6515
6520
6525
6530
6535
6540
6545
6550
6555
6560
6565
6570
6575
6580
6585
6590
6595
6600
6605
6610
6615
6620
6625
6630
6635
6640
6645
6650
6655
6660
6665
6670
6675
6680
6685
6690
6695
6700
6705
6710
6715
6720
6725
6730
6735
6740
6745
6750
6755
6760
6765
6770
6775
6780
6785
6790
6795
6800
6805
6810
6815
6820
6825
6830
6835
6840
6845
6850
6855
6860
6865
6870
6875
6880
6885
6890
6895
6900
6905
6910
6915
6920
6925
6930
6935
6940
6945
6950
6955
6960
6965
6970
6975
6980
6985
6990
6995
7000
7005
7010
7015
7020
7025
7030
7035
7040
7045
7050
7055
7060
7065
7070
7075
7080
7085
7090
7095
7100
7105
7110
7115
7120
7125
7130
7135
7140
7145
7150
7155
7160
7165
7170
7175
7180
7185
7190
7195
7200
7205
7210
7215
7220
7225
7230
7235
7240
7245
7250
7255
7260
7265
7270
7275
7280
7285
7290
7295
7300
7305
7310
7315
7320
7325
7330
7335
7340
7345
7350
7355
7360
7365
7370
7375
7380
7385
7390
7395
7400
7405
7410
7415
7420
7425
7430
7435
7440
7445
7450
7455
7460
7465
7470
7475
7480
7485
7490
7495
7500
7505
7510
7515
7520
7525
7530
7535
7540
7545
7550
7555
7560
7565
7570
7575
7580
7585
7590
7595
7600
7605
7610
7615
7620
7625
7630
7635
7640
7645
7650
7655
7660
7665
7670
7675
7680
7685
7690
7695
7700
7705
7710
7715
7720
7725
7730
7735
7740
7745
7750
7755
7760
7765
7770
7775
7780
7785
7790
7795
7800
7805
7810
7815
7820
7825
7830
7835
7840
7845
7850
7855
7860
7865
7870
7875
7880
7885
7890
7895
7900
7905
7910
7915
7920
7925
7930
7935
7940
7945
7950
7955
7960
7965
7970
7975
7980
7985
7990
7995
8000
8005
8010
8015
8020
8025
8030
8035
8040
8045
8050
8055
8060
8065
8070
8075
8080
8085
8090
8095
8100
8105
8110
8115
8120
8125
8130
8135
8140
8145
8150
8155
8160
8165
8170
8175
8180
8185
8190
8195
8200
8205
8210
8215
8220
8225
8230
8235
8240
8245
8250
8255
8260
8265
8270
8275
8280
8285
8290
8295
8300
8305
8310
8315
8320
8325
8330
8335
8340
8345
8350
8355
8360
8365
8370
8375
8380
8385
8390
8395
8400
8405
8410
8415
8420
8425
8430
8435
8440
8445
8450
8455
8460
8465
8470
8475
8480
8485
8490
8495
8500
8505
8510
8515
8520
8525
8530
8535
8540
8545
8550
8555
8560
8565
8570
8575
8580
8585
8590
8595
8600
8605
8610
8615
8620
8625
8630
8635
8640
8645
8650
8655
8660
8665
8670
8675
8680
8685
8690
8695
8700
8705
8710
8715
8720
8725
8730
8735
8740
8745
8750
8755
8760
8765
8770
8775
8780
8785
8790
8795
8800
8805
8810
8815
8820
8825
8830
8835
8840
8845
8850
8855
8860
8865
8870
8875
8880
8885
8890
8895
8900
8905
8910
8915
8920
8925
8930
8935
8940
8945
8950
8955
8960
8965
8970
8975
8980
8985
8990
8995
9000
9005
9010
9015
9020
9025
9030
9035
9040
9045
9050
9055
9060
9065
9070
9075
9080
9085
9090
9095
9100
9105
9110
9115
9120
9125
9130
9135
9140
9145
9150
9155
9160
9165
9170
9175
9180
9185
9190
9195
9200
9205
9210
9215
9220
9225
9230
9235
9240
9245
9250
9255
9260
9265
9270
9275
9280
9285
9290
9295
9300
9305
9310
9315
9320
9325
9330
9335
9340
9345
9350
9355
9360
9365
9370
9375
9380
9385
9390
9395
9400
9405
9410
9415
9420
9425
9430
9435
9440
9445
9450
9455
9460
9465
9470
9475
9480
9485
9490
9495
9500
9505
9510
9515
9520
9525
9530
9535
9540
9545
9550
9555
9560
9565
9570
9575
9580
9585
9590
9595
9600
9605
9610
9615
9620
9625
9630
9635
9640
9645
9650
9655
9660
9665
9670
9675
9680
9685
9690
9695
9700
9705
9710
9715
9720
9725
9730
9735
9740
9745
9750
9755
9760
9765
9770
9775
9780
9785
9790
9795
9800
9805
9810
9815
9820
9825
9830
9835
9840
9845
9850
9855
9860
9865
9870
9875
9880
9885
9890
9895
9900
9905
9910
9915
9920
9925
9930
9935
9940
9945
9950
9955
9960
9965
9970
9975
9980
9985
9990
9995
10000
10005
10010
10015
10020
10025
10030
10035
10040
10045
10050
10055
10060
10065
10070
10075
10080
10085
10090
10095
10100
10105
10110
10115
10120
10125
10130
10135
10140
10145
10150
10155
10160
10165
10170
10175
10180
10185
10190
10195
10200
10205
10210
10215
10220
10

turning the X and/or Y rotation thumbwheels 608, 610 rotate the scene accordingly around the point of rotation.

In a preferred embodiment, the right mouse button is reserved for the pop-up menus 616, 902. Pressing the right mouse button anywhere over an empty rendering area brings up the viewer pop-up menu 902. Pressing the right mouse button over an object brings up the object pop-up menu 616.

The viewer pop-up menu 902 allows the user to select the mode (such modes are described above), change viewer properties (set up preferences, e.g. background color), toggle on/off sets of objects, and/or access any other viewer features.

The object pop-up menu 616 allows the user to change an object's shape, color (material), and/or size, select the corresponding set of compounds, and/or define a neighborhood 3D area around the object (zoom feature, described below). In a preferred embodiment, all changes made to an object automatically apply to all other objects from the same set. The object's shape can be changed to one of the predefined basic shapes (e.g. dot, cube, sphere, cone). The object's material (color) is changed via a color dialog. The object's size is changed via a resize dialog. Any set of objects can be visible (toggled on) or hidden (toggled off). A toggle sets command brings up a list of sets defined for the current map 640. Clicking on a set in the list (highlighting/clearing) toggles the set off and on.

Invoking the zoom feature (via the pick neighbors command on the object pop-up menu 616, for example) creates a sphere 704 in the render area 614 (FIG. 7), which is centered on the object. The radius of the sphere 704 can be adjusted via a resize dialog 702 to select a desired neighborhood area around the object. All objects (and corresponding compounds) encompassed by the sphere 704 are then selected, displayed in a different map, added to a new or existing set, dragged to a target (described below), and/or viewed in a structure browser window 502.

The map viewer 112 is capable of maintaining an interactive selection of objects/compounds. All selected objects are visualized in the same shape, color, and/or size. In other words, selecting an object changes its shape, color, and/or size (e.g. to a purple cone), deselecting an object changes its shape, color and/or size back

to the original attributes. Executing the select set command from the object pop-up menu 616 selects the whole set of objects this object belongs to. Alternatively, an individual object can be selected or deselected by clicking a middle mouse button over an object. The interactive selection of objects can be converted to a set of compounds and displayed in a structure browser window 502. The current selection can be converted into a set of compounds by invoking the save selection command from a selection menu, and/or it can be cleared by executing the clear selection command from the selection menu.

9.3 *Interactivity of the Present Invention*

As should be apparent from the above, the present invention enables users to interact with the objects/compounds displayed in a compound visualization non-linear map. This interactivity provided by the present invention shall be further illustrated below.

9.3.1 *Map Viewer as Target*

According to the present invention, a user can select a plurality of compounds from some source, and then add those compounds to a new or an existing compound visualization non-linear map being displayed in a map window 600. In this instance, the map window 600 (or, equivalently in this context, the map viewer 112) is acting as a target for an interactive user activity.

This operation is conceptually shown in FIG. 14. A compound visualization non-linear map 1404 is being displayed in a map window 600. According to the present invention, the user can select compounds from a structure browser window 502, and then add those selected compounds (through, for example, well known drag and drop operations) to the compound visualization non-linear map 1404. Similarly, the user can select compounds from a compound database 122, or from a MS (mass

spectrometry) viewer 1402, and then add those compounds to the compound visualization non-linear map 1404.

According to an embodiment of the invention, new compounds are added to an existing compound visualization non-linear map by incremental refinement of the compound visualization non-linear map. Such incremental refinement is described above.

9.3.2 Map Viewer as Source

According to the present invention, a user can select a plurality of compounds from a map window 600, and then have those compounds processed by a target. In this instance, the map window 600 (or, equivalently in this context, the map viewer 112) is acting as a source for an interactive user activity.

This operation is conceptually shown in FIG. 13. A user selects one or more compounds from the compound visualization non-linear map being displayed in the map window 600, and then drags and drops the selected compounds to a target. The described action is interpreted as a submission of the corresponding chemical structure(s) to the receiving target for processing. The receiving object can be anything that can handle a chemical structure: another map viewer 112, a structure viewer 110, a (molecular) spreadsheet 136, a database 120, an experiment planner 140, an active site docker 144, an NMR widget 130, an MS widget 134, a QSAR model 138, a property prediction program 142, or any other suitable process. For example, dragging and dropping a compound onto an NMR widget would display this compound's NMR spectrum, either an experimental or a predicted one.

The experiment planner is described in pending U.S. Patent Application titled "SYSTEM, METHOD AND COMPUTER PROGRAM PRODUCT FOR IDENTIFYING CHEMICAL COMPOUNDS HAVING DESIRED PROPERTIES," Atty. Docket No. 1503.0200001, herein incorporated by reference in its entirety.

The drag and drop concept described above provides a powerful enhancement of a 3D mapping and visualization of compound collections and libraries. Any

conceivable information about a set of chemical compounds can thus be easily accessed from the compound visualization non-linear map. For example, a map of compounds capable of binding to an active site of a given enzyme or receptor would benefit from the possibility to visualize how compounds from the different areas of the map bind to that enzyme or receptor.

9.4 Multiple Maps

According to the present invention, it is possible to create multiple visual maps for any given set of collections and/or libraries of chemical compounds. Multiple visual maps can be based on the same and/or different non-linear maps. Visual maps based on the same non-linear map can display different subsets of compounds and/or present different views of the same set of compounds (e.g. one visual map can display an XY plane view and another visual map can display an orthogonal, YZ plane view). Visual maps based on different non-linear maps can visualize the same set of compounds on different projections, for example, maps derived from different similarity relations between these compounds.

If a compound is mapped on multiple visual maps, the visual objects representing the compound on the different maps can be crosslinked. Crosslinking means that any modifications made to a visual object in one of the visual maps will be automatically reflected into the other visual maps. For example, if an object is selected on one of the visual maps, it will be displayed as selected on the other visual maps as well. In fact, all objects on all maps can be crosslinked provided that they represent the same chemical compounds. Multiple visual maps can be also crosslinked in a way that mapping any additional compounds onto one of the visual maps will automatically map the same compounds onto the crosslinked maps.

10. Examples

The present invention is useful for visualizing and interactively processing any chemical entities including but not limited to small molecules, polymers, peptides, proteins, etc. It may also be used to display different similarity relationships between these compounds.

5 The present invention has been described above with the aid of functional building blocks illustrating the performance of specified functions and relationships thereof. The boundaries of these functional building blocks have been arbitrarily defined herein for the convenience of the description. Alternate boundaries can be defined so long as the specified functions and relationships thereof are appropriately performed. Any such alternate boundaries are thus within the scope and spirit of the claimed invention. These functional building blocks may be implemented by discrete components, application specific integrated circuits, processors executing appropriate software and the like or any combination thereof. It is well within the scope of one skilled in the relevant art(s) to develop the appropriate circuitry and /or software to implement these functional building blocks.

10 While various embodiments of the present invention have been described above, it should be understood that they have been presented by way of example only, and not limitation. Thus, the breadth and scope of the present invention should not be limited by any of the above-described exemplary embodiments, but should be defined only in accordance with the following claims and their equivalents.